# International Metadata Standards and Enterprise Data Quality Metadata Systems

IN027 - December 15, 2016

Ted Habermann

Director of Earth Science

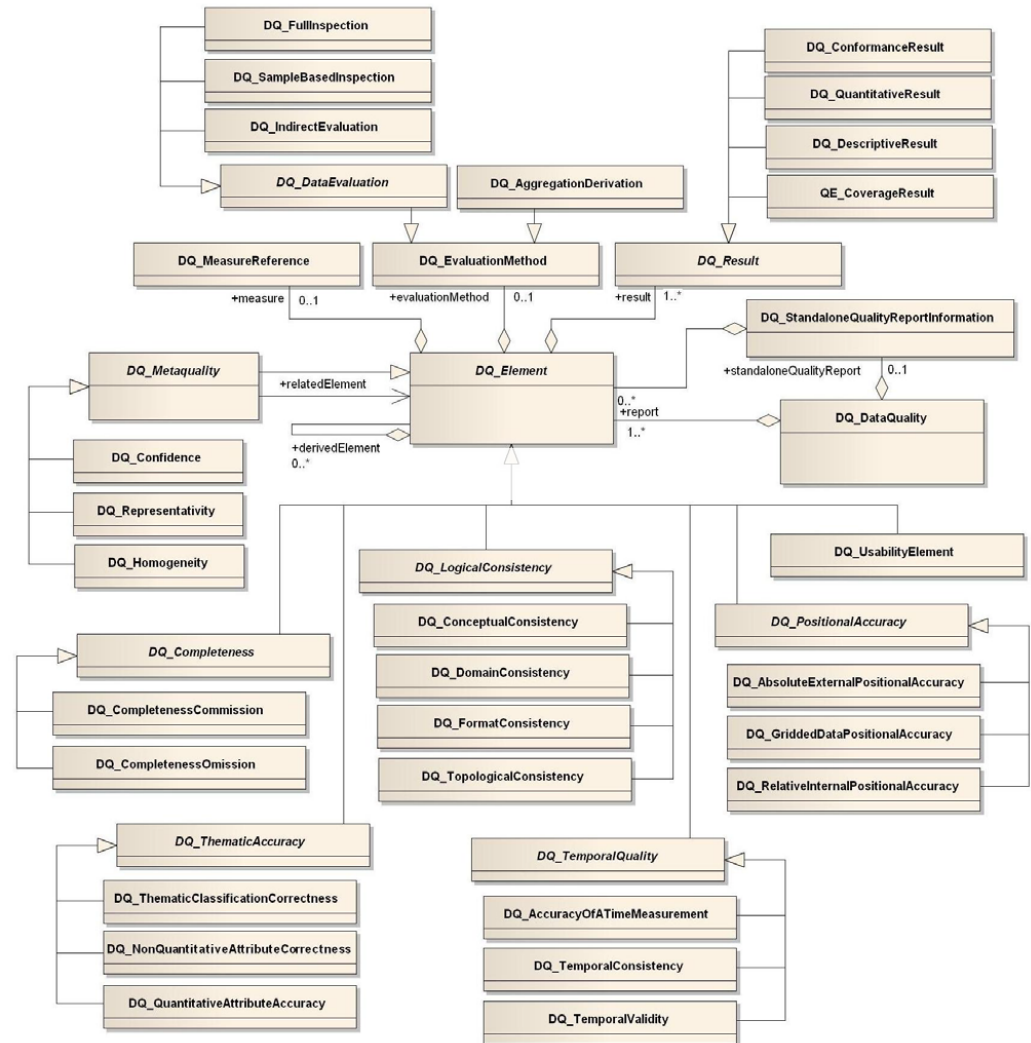The HDF Group

thabermann@hdfgroup.org

# The Big Picture

ISO 19157 is a conceptual model of data quality metadata that was recently approved as an international standard. It combines three older standards into a unified model for describing data quality.

Many of the principle elements of this conceptual model are abstract, and can be implemented in several ways.
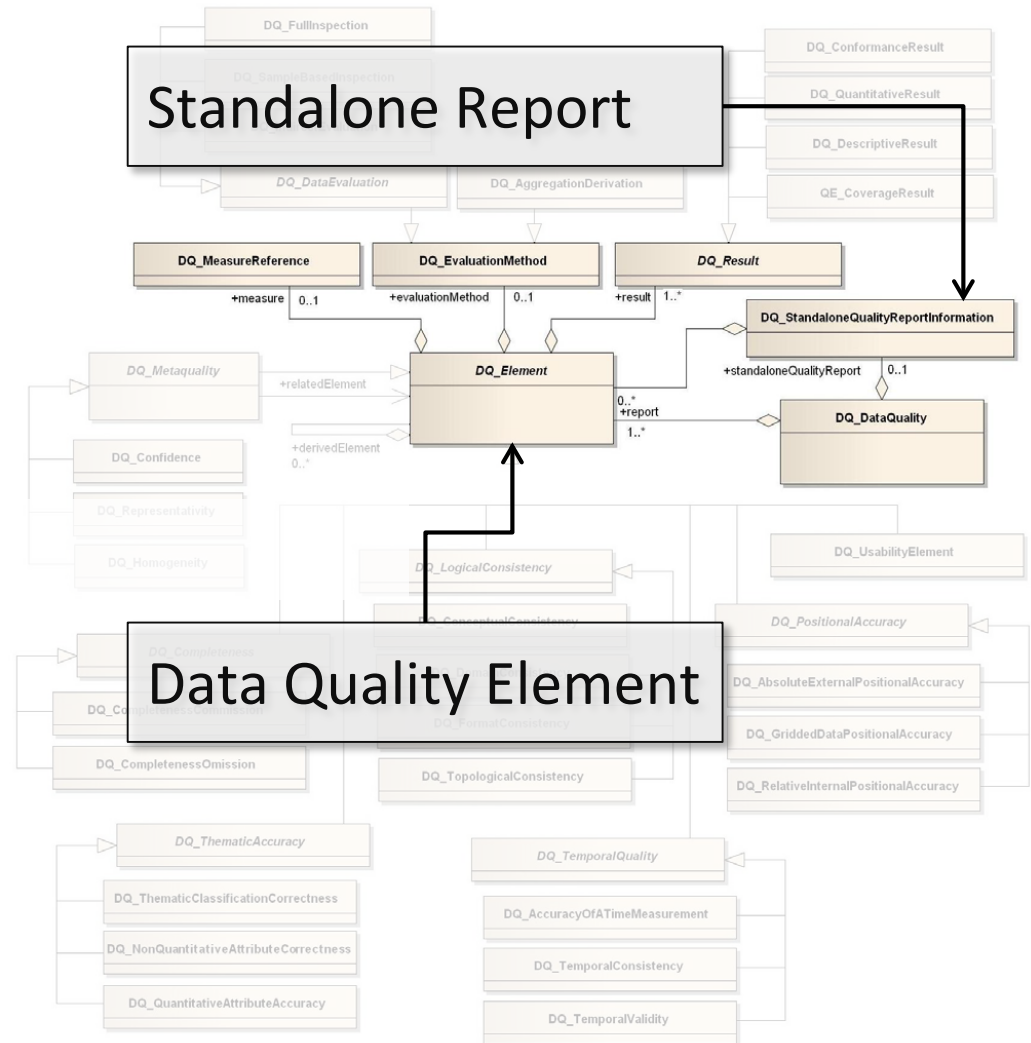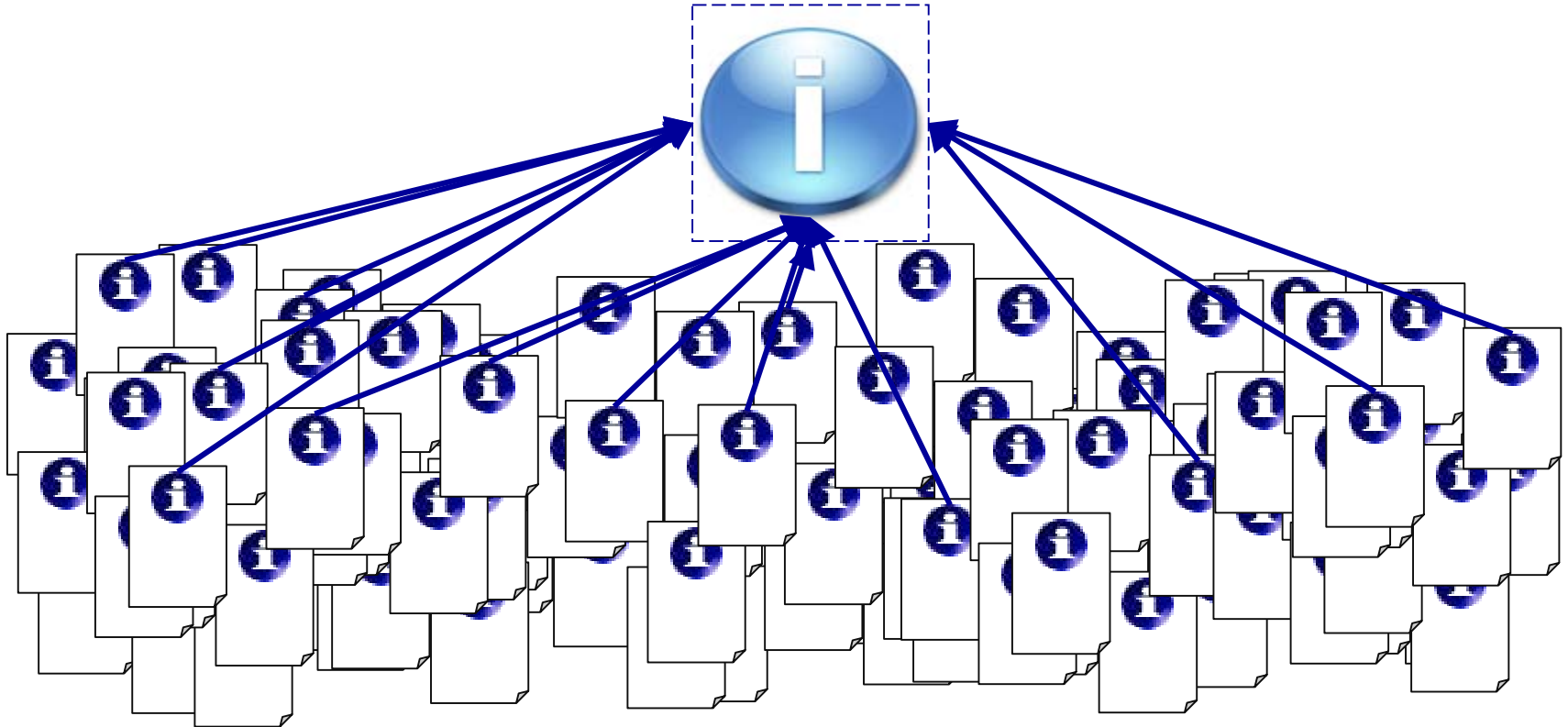
# The Big Picture

ISO 19157 is a conceptual model of data quality metadata that was recently approved as an international standard. It combines three older standards into a unified model for describing data quality.

Many of the principle elements of this conceptual model are abstract, they can be implemented in several ways.

When only the abstract concepts are considered, the model is very simple.



**EOSDIS**

# Enterprise Systems?

# Stand Alone Quality Reports

*"There are papers and web pages that describe the quality of my data."*

Papers and reports that describe data quality are StandAloneReports. Metadata can include brief descriptions of the results (abstracts) and references to any number of these (citations).

| DQ_StandaloneQualityReportInformation |
|---|
| + abstract : CharacterString<br>+ reportReference: CI_Citation |

Abstract: The fire training-set may also have been biased against savanna and savanna woodland fires since their detection is more difficult than in humid, forest environments with cool background temperatures [Malingreau, 1990]. There may, therefore, be an under-sampling of warmer background environments.
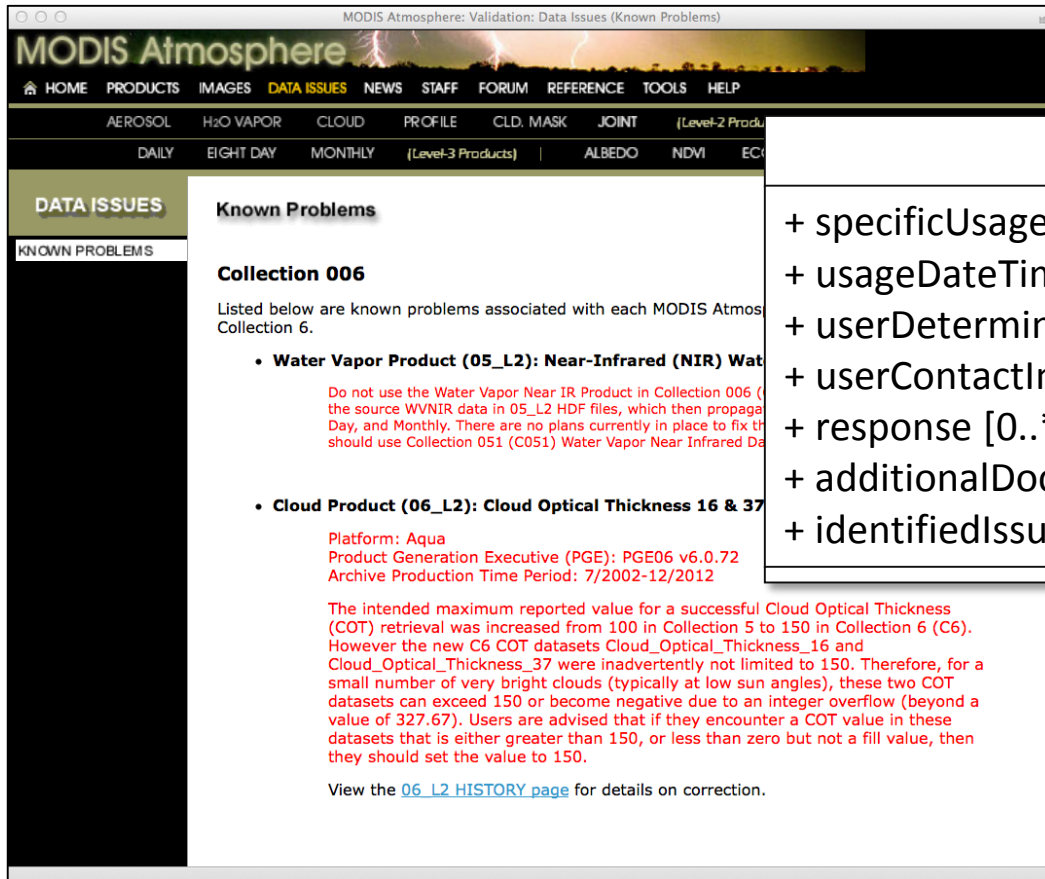
DOI

Citation: Malingreau J.P, 1990, The contribution of remote sensing to the global monitoring of fires in tropical and subtropical ecosystems. In: *Fire in Tropical Biota*, (J.G. Goldammer , editor), Springer Verlag , Berlin: 337-370.

# Data Usage (19115-1)

*"Users increase our understanding of data quality. We need to keep them in the loop."*



MD_Usage

+ specificUsage : CharacterString
+ usageDateTime [0..1] : DateTime
+ userDeterminedLimitations [0..1] : CharacterString
+ userContactInfo [1..*] : CI_ResponsibleParty
+ response [0..*] : CharacterString
+ additionalDocumentation [0..*] : CI_Citation
+ identifiedIssues [0..1] : CI_Citation

DOI

# What is a Data Quality Element?



QA_PercentMissingData

Data Quality Element

Measure

Result

15%

Method

$$\frac{\text{Number of Pixels with Missing Flags}}{\text{Total Number of Pixels}}$$

EOSDIS

# What Are Quality Measures?

*"My metadata already include data quality measures ."*

NASA EOSDIS metadata includes two types of quality measures.

**QA_Stats**

**QA_Flags**

### 4.7  MEASURED PARAMETERS

Measured parameters are associated only at the granule level only and are important search criteria for granules. For some providers, the value of certain measured parameters determi

Measured parameters contain the name of the geophysical parameter associated quality flags and quality status. The quality status contains parameters used to set these measures are not preset and will be dete measures can occur many times either for the granule as a whole or fo contain the science, operational and automatic quality flags that indicat specific parameter values within a granule.

A measured parameter is uniquely identified by its **ParameterName** ele

- **QAStats** – The name of the geophysical parameter expressed flags and quality status.
  - o **QAPercentMissingData** - Granule level % missing da individual parameters within a granule.
  - o **QAPercentOutOfBoundsData** – Granule level % out repeated for individual parameters within a granule.

*ECHO 10.0 Data Partner's User Guide's Data Partner's User Guide*

Version: 10.7
March 2010

- o **QAPercentInterpolatedData** – Granule level % interpolated data. This attribute can be repeated for individual parameters within a granule.
- o **QAPercentCloudCover** – This attribute is used to characterize the cloud cover amount of a granule. This attribute may be repeated for individual parameters within a granule. (Note - there may be more than one way to define a cloud or it's effects within a product containing several parameters; i.e. this attribute may be parameter specific)

- **QAFlags** – The name of the geophysical parameter expressed in the data as well as associated quality flags and quality status.
  - o **AutomaticQualityFlag** – The granule level flag applying generally to the granule and specifically to parameters the granule level. When applied to parameter, the flag refers to the quality of that parameter for the granule (as applicable). The parameters determining whether the flag is set are defined by the developer and documented in the Quality Flag Explanation.
  - o **AutomaticQualityFlagExplanation** – A text explanation of the criteria used to set automatic quality flag, including thresholds or other criteria.
  - o **OperationalQualityFlag** – The granule level flag applying both generally to a granule and specifically to parameters at the granule level. When applied to parameter, the flag refers to the quality of that parameter for the granule (as applicable). The parameters determining whether the flag is set are defined by the developers and documented in the Operational Quality Flag Explanation.
  - o **OperationalQualityFlagExplanation** – A text explanation of the criteria used to set operational quality flag; including thresholds or other criteria.
  - o **ScienceQualityFlag** – Granule level flag applying to a granule, and specifically to parameters. When applied to parameter, the flag refers to the quality of that parameter for the granule (as applicable). The parameters determining whether the flag is set are defined by the developers and documented in the Science Quality Flag Explanation.
  - o **ScienceQualityFlagExplanation** – A text explanation of the criteria used to set science quality flag; including thresholds or other criteria.

**EOSDIS**

# What Are Quality Measures?

*"I use consistent Quality Measures across many products."*

QA_Stats

### 4.7 MEASURED PARAMETERS

Measured parameters are associated only at the granule ... For some providers, the value of certain measured parameters determines the visibility of the granule.

Measured parameters contain the name of the geophysical parameter expressed in the data as well as associated quality flags and quality status. The quality status contains measures of quality for the granule. The parameters used to set these measures are not preset and will be determined by the data producer. Each set of measures can occur many times either for the granule as a whole or for individual parameters. The quality flags contain the science, operational and automatic quality flags that indicate the overall quality assurance levels of specific parameter values within a granule.

A measured parameter is uniquely identified by its **ParameterName** element, and has the following information:

- **QAStats** – The name of the geophysical parameter expressed in the data as well as associated quality flags and quality status.
  - **QAPercentMissingData** - Granule level % missing data. This attribute can be repeated for individual parameters within a granule.
  - **QAPercentOutOfBoundsData** – Granule level % out of bounds data. This attribute can be repeated for individual parameters within a granule.

*ECHO 10.0 Data Partner's User Guide's Data Partner's User Guide*      Page 57

Version: 10.7
March 2010

**QAStats** – **Standard measures for all products**
**QAPercentMissingData** - Granule level % missing data. This attribute can be repeated for individual parameters within a granule.
**QAPercentOutOfBoundsData** – Granule level % out of bounds data. This attribute can be repeated for individual parameters within a granule.
**QAPercentInterpolatedData** – Granule level % interpolated data. This attribute can be repeated for individual parameters within a granule.
**QAPercentCloudCover** – This attribute is used to characterize the cloud cover amount of a granule. This attribute may be repeated for individual parameters within a granule. (Note - there may be more than one way to define a cloud or it's effects within a product containing several parameters; i.e. this attribute may be parameter specific)

![EOSDIS]

# What Are Quality Measures?

*"I use consistent types of Quality Measure across many products."*

**QA_Flags**

**QAFlags** – **Classes of quality measures with product specific implementations**

**AutomaticQualityFlag** – The granule level flag applying generally to the granule and specifically to parameters the granule level. When applied to parameter, the flag refers to the quality of that parameter for the granule (as applicable). The parameters determining whether the flag is set are defined by the developer and documented in the Quality Flag Explanation.

**AutomaticQualityFlagExplanation** – A text explanation of the criteria used to set automatic quality flag, including thresholds or other criteria.

**OperationalQualityFlag** – The granule level flag applying both generally to a granule and specifically to parameters at the granule level. When applied to parameter, the flag refers to the quality of that parameter for the granule (as applicable). The parameters determining whether the flag is set are defined by the developers and documented in the Operational Quality Flag Explanation.

**OperationalQualityFlagExplanation** – A text explanation of the criteria used to set operational quality flag; including thresholds or other criteria.

**ScienceQualityFlag** – Granule level flag applying to a granule, and specifically to parameters. When applied to parameter, the flag refers to the quality of that parameter for the granule (as applicable). The parameters determining whether the flag is set are defined by the developers and documented in the Science Quality Flag Explanation.

**ScienceQualityFlagExplanation** – A text explanation of the criteria used to set science quality flag; including thresholds or other criteria.
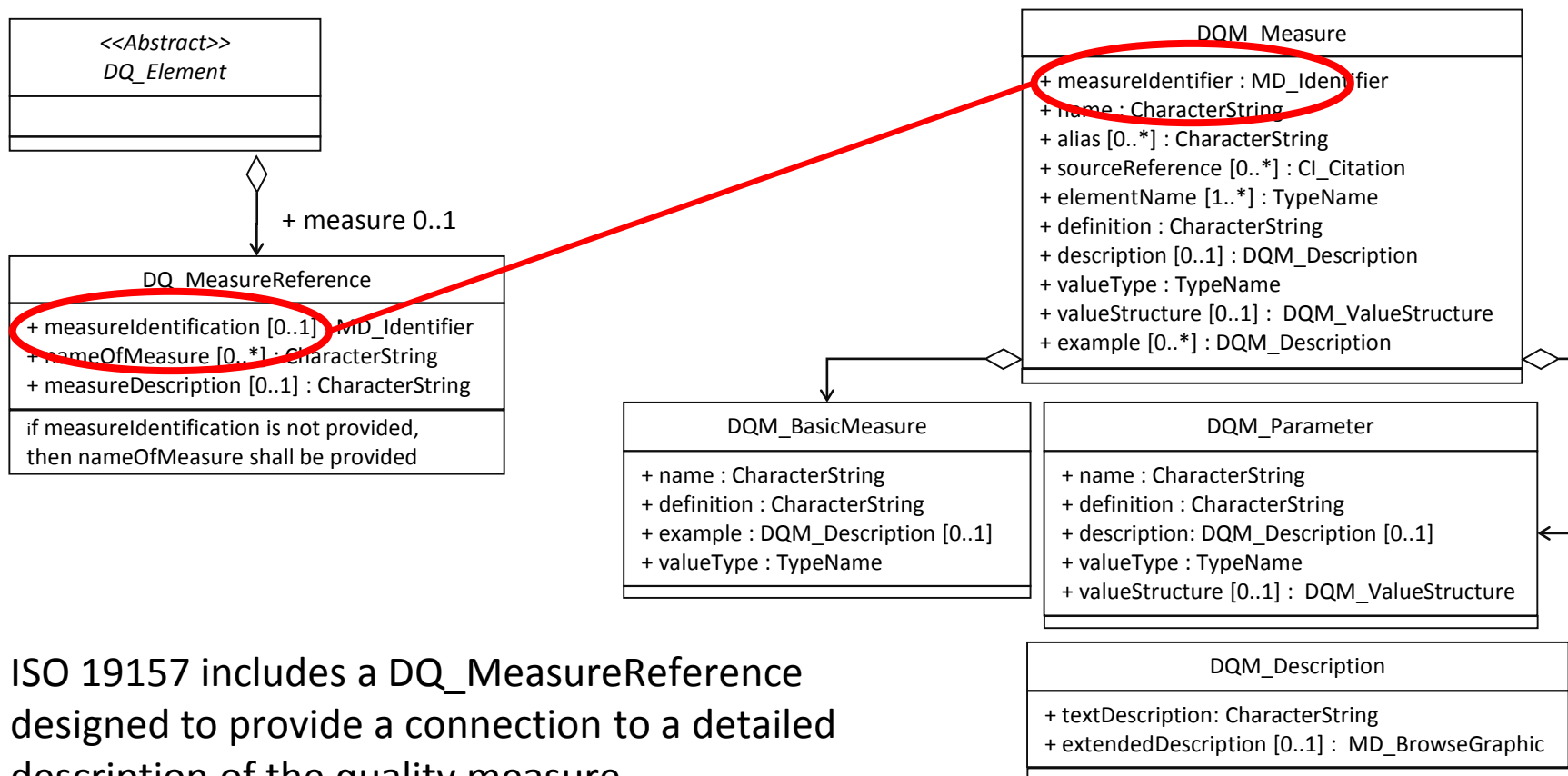
# Data Quality Measures

*"My data quality measures are consistently described in a database ."*



**DQ_Element**
<<Abstract>>

**DQ_MeasureReference**
+ measureIdentification [0..1] : MD_Identifier
+ nameOfMeasure [0..*] : CharacterString
+ measureDescription [0..1] : CharacterString

if measureIdentification is not provided,
then nameOfMeasure shall be provided

+ measure 0..1

**DQM_Measure**
+ measureIdentifier : MD_Identifier
+ name : CharacterString
+ alias [0..*] : CharacterString
+ sourceReference [0..*] : CI_Citation
+ elementName [1..*] : TypeName
+ definition : CharacterString
+ description [0..1] : DQM_Description
+ valueType : TypeName
+ valueStructure [0..1] : DQM_ValueStructure
+ example [0..*] : DQM_Description

**DQM_BasicMeasure**
+ name : CharacterString
+ definition : CharacterString
+ example : DQM_Description [0..1]
+ valueType : TypeName

**DQM_Parameter**
+ name : CharacterString
+ definition : CharacterString
+ description: DQM_Description [0..1]
+ valueType : TypeName
+ valueStructure [0..1] : DQM_ValueStructure

**DQM_Description**
+ textDescription: CharacterString
+ extendedDescription [0..1] : MD_BrowseGraphic

ISO 19157 includes a DQ_MeasureReference designed to provide a connection to a detailed description of the quality measure.

**EOSDIS**

# Data Quality Measures

*"I need to clearly and consistently explain how I measure quality."*

The ISO model for quality measures
includes identifiers, definitions,
descriptions, references and illustrations.

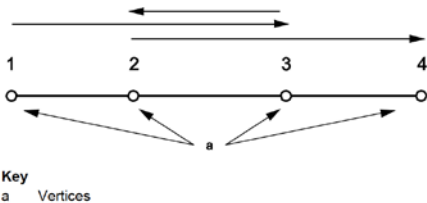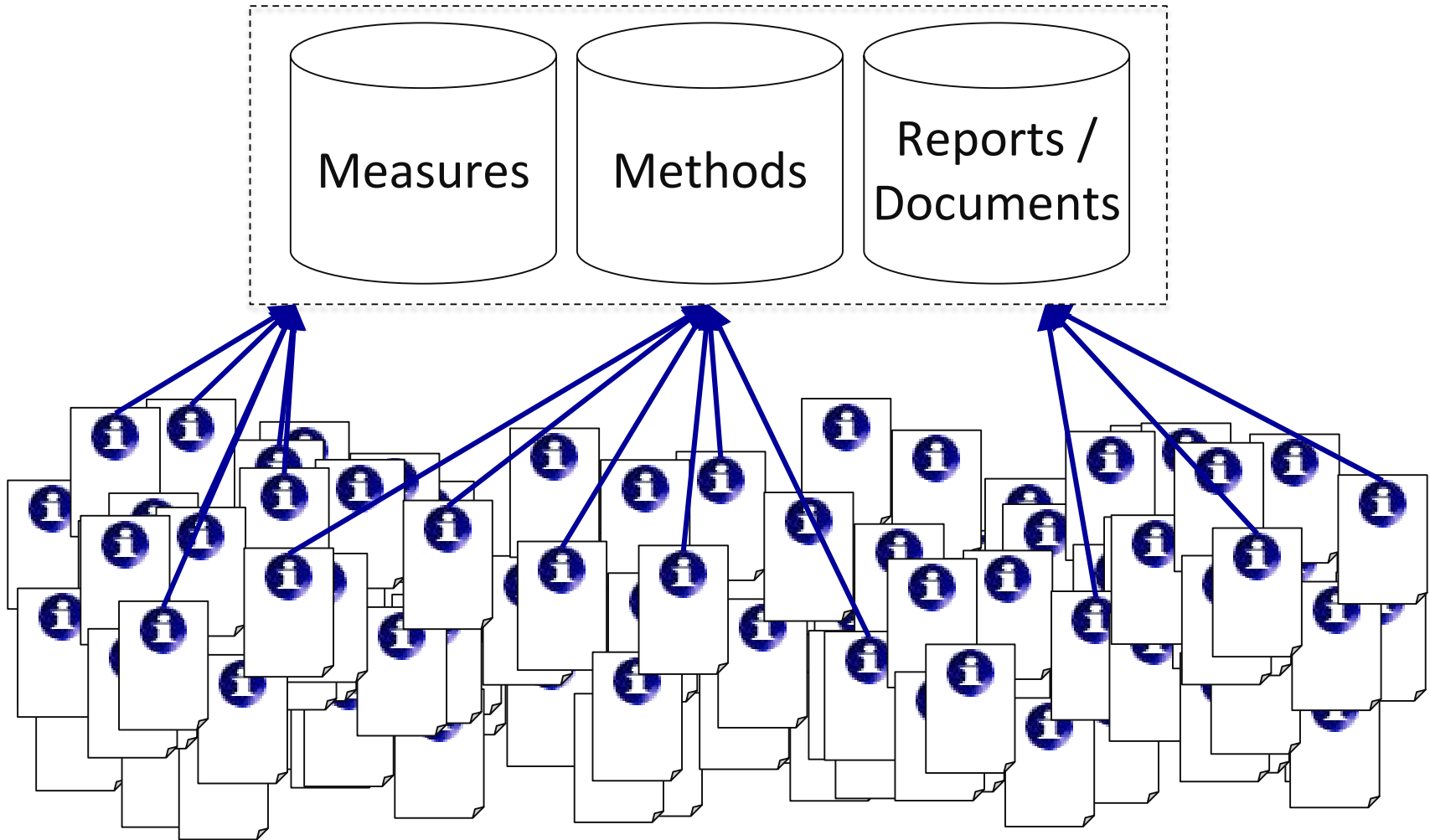**Table D.28 — Number of invalid self-overlap errors**

| Line | Component | Description |
|---|---|---|
| 1 | Name | number of invalid self-overlap errors |
| 2 | Alias | kickbacks |
| 3 | Element name | topological consistency |
| 4 | Basic measure | error count |
| 5 | Definition | count of all items in the data that illegally self overlap |
| 6 | Description | – |
| 7 | Parameter | – |
| 8 | Value type | Integer |
| 9 | Value structure | – |
| 10 | Source reference | – |
| 11 | Example | (illustration) Key: a Vertices |
| 12 | Identifier | 27 |

**Table D.31 — Mean value of positional uncertainties excluding outliers**

| Line | Component | Description |
|---|---|---|
| 1 | Name | mean value of positional uncertainties excluding outliers (2D) |
| 2 | Alias | – |
| 3 | Element name | absolute or external accuracy |
| 4 | Basic measure | not applicable |
| 5 | Definition | for a set of points where the distance does not exceed a defined threshold, the arithmetical average of distances between their measured positions and what is considered as the corresponding true positions |
| 6 | Description | For a number of points ($N$), the measured positions are given as $x_{mi}$, $y_{mi}$ and $z_{mi}$ coordinates depending on the dimension in which the position of the point is measured. A corresponding set of coordinates, $x_{ti}$, $y_{ti}$ and $z_{ti}$, are considered to represent the true positions. All positional uncertainties above a defined threshold $e_{max}$ are then removed from the set. The positional uncertainties are calculated as $$e_i^{\cdot} = \begin{cases} e_i, & if \quad e_i \leq e_{max} \\ 0, & if \quad e_i > e_{max} \end{cases}$$ The calculation of $e_i$ is given by the data quality measure "mean value of positional uncertainties" in one, two and three dimensions. For the remaining number of errors ($N_R$), the mean of the horizontal absolute positions is calculated as $$\bar{a}_{\text{excluding outliers}} = \frac{1}{N_R} \sum_{i=1}^{N} e_i^{\cdot}$$ A criterion for the establishing of correspondence should also be stated (e.g. allowing for correspondence to the closest position, correspondence on vertices or along lines). The criteria for finding the corresponding points shall be reported with the data quality evaluation result. |
| 7 | Parameter | Name: $e_{max}$ Definition: is the threshold for accepted positional uncertainties Value type: Number |
| 8 | Value type | Measure |
| 9 | Value structure | – |
| 10 | Source reference | – |
| 11 | Example | – |
| 12 | Identifier | 29 |

**EOSDIS**

12

# Enterprise Systems?

# Summary



"There are papers and web pages that describe the quality of my data."

"My metadata currently includes descriptions of the quality of my data."

"Users increase our understanding of data quality. We need to keep them in the loop."

"My data quality information exists in databases or web services."

"I use consistent types of Quality Measure across many products."

"The quality of my data vary in time and space and different parameters have different quality measures and results."

"I use consistent Quality Measures across many products."

"I need to clearly and consistently explain how I measure quality."

# Acknowledgements



This work was partially supported by contract number NNG15HZ39C from NASA.

Any opinions, findings, conclusions, or recommendations expressed in this material are those of the author and do not necessarily reflect the views of NASA or The HDF Group.